

WE CLAIM:

1. An interconnection controller, comprising:

an intra-cluster interface configured for coupling with intra-cluster links to a plurality
5 of local nodes arranged in a point-to-point architecture in a local cluster, the local nodes
including local processors;

an inter-cluster interface configured for coupling with an inter-cluster link to a non-
local interconnection controller in a non-local cluster;

encapsulation logic configured to receive intra-cluster packets from the local nodes
10 via the intra-cluster links and to encapsulate the intra-cluster packets as inter-cluster packets
for transmission on the inter-cluster link; and

a module comprising a remote transmission buffer, the module configured to:

receive inter-cluster packets from the encapsulation logic;

store inter-cluster packets in the remote transmission buffer;

15 forward inter-cluster packets for transmission on the inter-cluster link;

determine when the remote transmission buffer is empty;

generate a special packet for transmission on the inter-cluster link when the
buffer is empty; and

forward the special packet for transmission on the inter-cluster link without
20 storing the special packet in the remote transmission buffer.

2. The interconnection controller of claim 1, wherein the special packet
comprises a control character.

3. The interconnection controller of claim 1, the module further comprising a reception buffer, the module being further configured to:

receive a special packet from the inter-cluster link; and

drop the special packet without storing the special packet in the reception buffer.

5

4. The interconnection controller of claim 1, wherein the module is configured to determine when the transmission buffer is empty by inspecting a single buffer space of the transmission buffer.

10

5. The interconnection controller of claim 1, wherein the transmission buffer is an asynchronous buffer that is configured to receive inter-cluster packets from the encapsulation logic at a first clock speed and forwards the inter-cluster packets at a second clock speed.

15

6. The interconnection controller of claim 1, wherein the module is configured to initialize the inter-cluster link and use information obtained during an initialization process to perform de-skewing operations on packets received on the inter-cluster link.

20

7. The interconnection controller of claim 1, further comprising a local transmitter configured to:

receive intra-cluster packets from the from the local nodes via the intra-cluster links;

store intra-cluster packets in a local transmission buffer;

forward intra-cluster packets for transmission on the intra-cluster links;

25

determine when there are no valid intra-cluster packets to transmit;

generate NOP packets when there are no valid intra-cluster packets to transmit;

forward the NOP packets to the local transmission buffer; and
transmit the NOP packets on the intra-cluster links.

5

8. The interconnection controller of claim 1, further comprising a serializer for serializing inter-cluster packets and for performing bit conversion of inter-cluster packets.

9. An integrated circuit comprising the interconnection controller of claim 1.

10

10. A set of semiconductor processing masks representative of at least a portion of the interconnection controller of claim 1.

11. At least one computer-readable medium having data structures stored therein
15 representative of the interconnection controller of claim 1.

12. The interconnection controller of claim 8, wherein the bit conversion encodes clock data in the inter-cluster packets.

20 13. The interconnection controller of claim 8, wherein the bit conversion comprises 8b/10b conversion.

14. The integrated circuit of claim 9, wherein the integrated circuit comprises an application-specific integrated circuit.

25

15. The at least one computer-readable medium of claim 11, wherein the data structures comprise a simulatable representation of the interconnection controller.

16. The at least one computer-readable medium of claim 11, wherein the data structures comprise a code description of the interconnection controller.

17. The at least one computer-readable medium of claim 15, wherein the simulatable representation comprises a netlist.

18. The at least one computer-readable medium of claim 16, wherein the code description corresponds to a hardware description language.

19. A computer system, comprising:

a first cluster including a first plurality of processors and a first interconnection controller, the first plurality of processors and the first interconnection controller interconnected by first point-to-point intra-cluster links, the first interconnection controller comprising:

encapsulation logic configured to receive intra-cluster packets from the first plurality of processors via the first point-to-point intra-cluster links and to encapsulate the intra-cluster packets as high-speed link packets for transmission on an inter-cluster link; and

a first module comprising a transmission buffer, the first module configured to:

receive high-speed link packets from the encapsulation logic;

store high-speed link packets in a transmission buffer;

forward high-speed link packets for transmission on the inter-cluster link;
determine an empty condition indicating that the transmission buffer is
empty;
generate a special packet responsive to the empty condition; and
5 forward the special packet to the inter-cluster link without storing the special
packet in the transmission buffer.

20. The computer system of claim 19, further comprising:

a second cluster including a second plurality of processors and a second
10 interconnection controller, the second plurality of processors and the second interconnection
controller interconnected by second point-to-point intra-cluster links, the second
interconnection controller comprising a second module configured to:
receive high-speed link packets from the inter-cluster link;
store the high-speed link packets in a reception buffer,
15 receive the special packet; and
drop the special packet without storing the special packet in the reception buffer.

21. The computer system of claim 19, wherein the special packet comprises a
control character.

20

22. A computer system comprising a plurality of processor clusters
interconnected by a plurality of point-to-point inter-cluster links, each processor cluster
comprising nodes including a plurality of local processors and an interconnection controller
interconnected by a plurality of point-to-point intra-cluster links, communications within a
25 cluster being made via an intra-cluster protocol that uses intra-cluster packets, wherein the

interconnection controller in each cluster is operable to map locally-generated communications directed to others of the clusters to the point-to-point inter-cluster links and to map remotely-generated communications directed to the local nodes to the point-to-point intra-cluster links, communications between clusters being made via an inter-cluster protocol that uses inter-cluster packets, an inter-cluster packet encapsulating at least one intra-cluster packet, each interconnection controller configured to generate and transmit a special packet on an inter-cluster link when the interconnection controller has no valid inter-cluster packets to send, the special packet not being stored in a transmission buffer prior to being transmitted on the inter-cluster link.

23. The computer system of claim 22, wherein the special packet comprises a control character.

24. The computer system of claim 22, each interconnection controller being further configured to receive a special packet, but not to store the special packet in a reception buffer for storing valid inter-cluster packets.

25. The computer system of claim 22, wherein the transmission buffer is an asynchronous buffer that receives inter-cluster packets at a first clock speed and forwards the inter-cluster packets at a second clock speed for transmission on the inter-cluster link.

26. The computer system of claim 22, wherein each interconnection controller is further configured to generate a NOP packet when the interconnection controller has no valid intra-cluster packets to send.

27. The computer system of claim 26, wherein each interconnection controller is further configured to forward the NOP packet to a local transmission buffer to await transmission on an intra-cluster link.

5 28. A computer-implemented method for decreasing latency in a computer system comprising a plurality of clusters, each cluster including a plurality of local nodes and an interconnection controller interconnected by point-to-point intra-cluster links, communications between the local nodes and the interconnection controller made via an intra-cluster protocol using intra-cluster packets, the interconnection controller of each
10 cluster interconnected by inter-cluster links with the interconnection controller of other clusters, the computer-implemented method comprising:

forming inter-cluster packets by encapsulating intra-cluster packets;

storing the inter-cluster packets in a remote transmission buffer of a first interconnection controller;

15 transmitting the inter-cluster packets to a second interconnection controller;

storing received inter-cluster packets in a reception buffer of the second interconnection controller;

determining that the remote transmission buffer is empty;

generating a control character in response to a determination that the remote
20 transmission buffer is empty;

transmitting the control character to the second interconnection controller; and

dropping the control character without storing the control character in the reception buffer.

29. The computer-implemented method of claim 28, further comprising:
performing an initialization sequence that establishes a characteristic skew
pattern between data lanes of the inter-cluster link;
encoding clock data in each symbol transmitted on the inter-cluster link;
5 recovering clock data from each symbol received on the inter-cluster link; and
applying the characteristic skew pattern to correct for skew between data
lanes of the inter-cluster link.

30. The computer-implemented method of claim 28, further comprising:
10 determining that there are no valid intra-cluster packets for transmission on
the point-to-point intra-cluster links;
generating NOP packets;
storing the NOP packets in a local transmission buffer; and
transmitting the NOP packets to local nodes on the point-to-point intra-cluster
15 links.